# SRE & Happiness

Denis Ćutić
*Senior Site Reliability Engineer @ Infobip*

# Let me tell you a story...
# Story #1

29.04.2018

:(

# incident-management ⌄

**Denis Cutic**  10:38

what's happening?

# What happened?

**Denis Cutic** 13:22
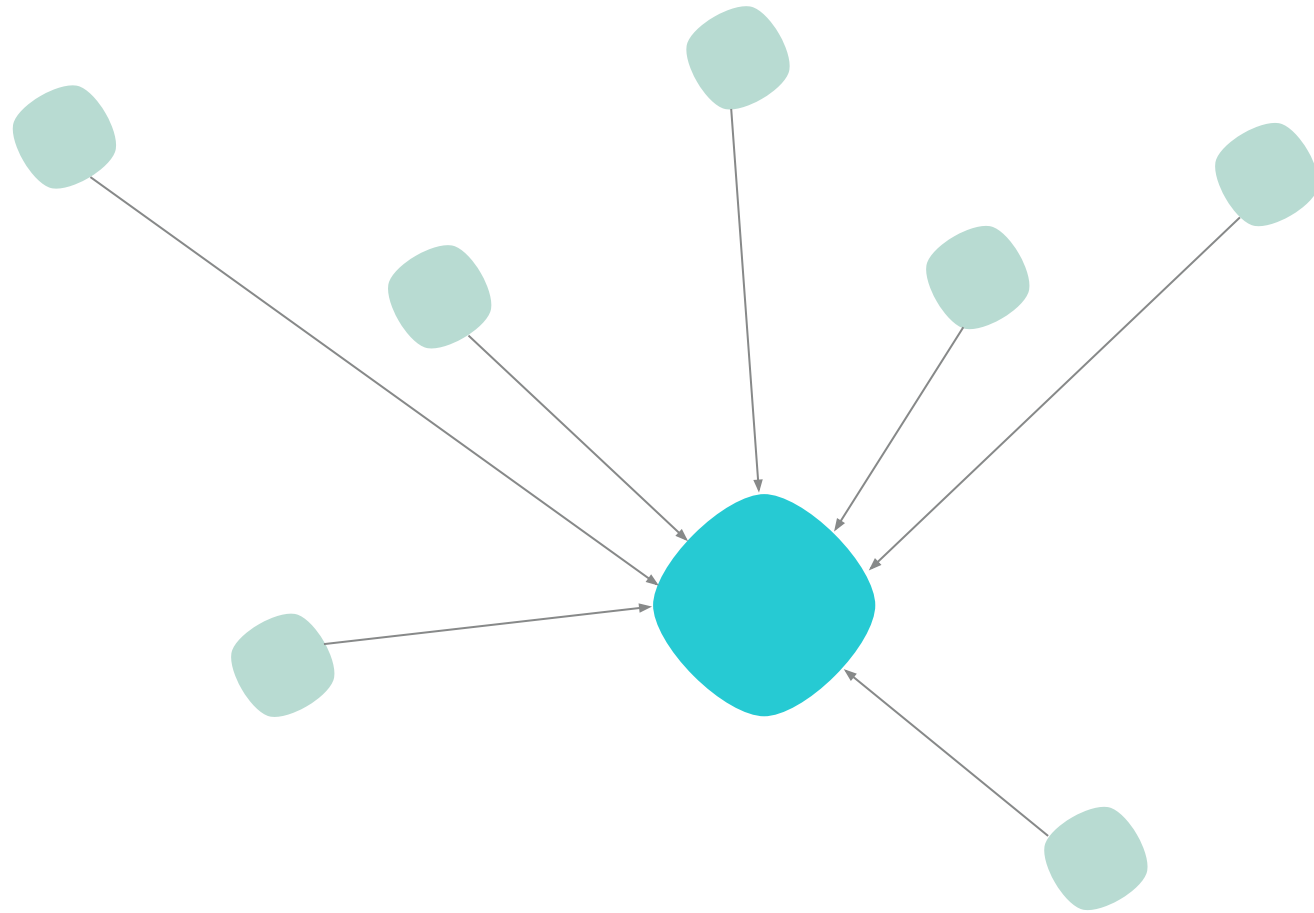
Dok mi se laptop hladi u frizideru mozes bacit oko jesu li

hope != strategy
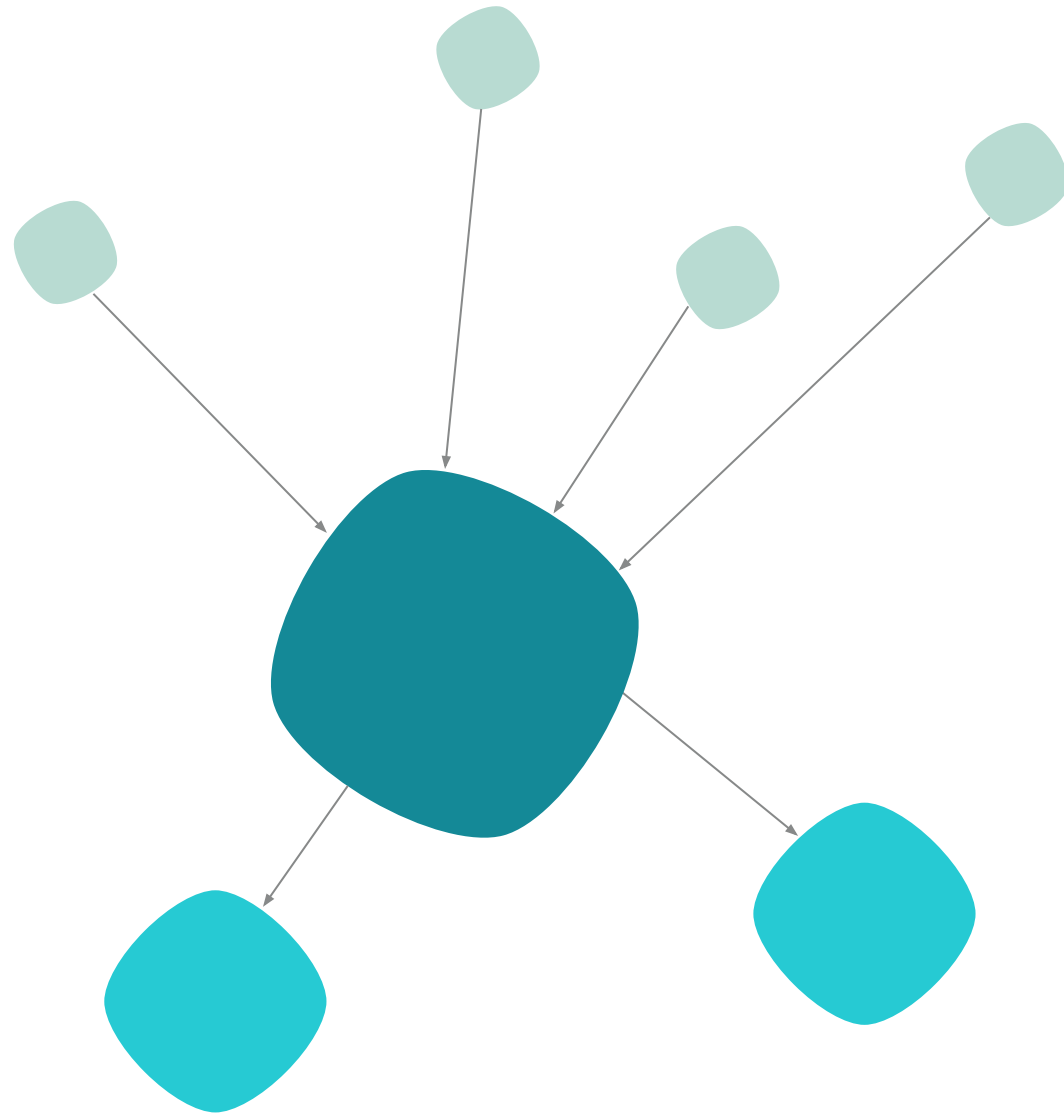
# My journey
# Story #2

Auth services
Account/User management services

# SMS API service

# Generic API service

# API infrastructure

# Survival?

SRE

**Infobip's journey**
**Story #3 - Prequel**

# DevOps VS SRE

# Infobip's journey
# Story #3

~~DevOps~~
Startup

**Reliability**

# Beginner level

- Reactive
- Hopefully someone will catch failures
- Hopefully it will be someone inside the company
- Hopefully someone will notice some patterns
- Hopefully we will meet the agreed SLA

Sys Admins + DevOps + Ops

**Reliability**

# Advanced level

- Reactive
- Observability tools in place
- Service-level monitoring from client PoV
- Promoted through company culture
- Support process in place
- Incident management process
- Starting with post-incident reviews

Sys Admins + DevOps + Ops +
QA + SecOps + SRE

**Reliability**

# Shipping stuff to space level

Not there yet...

# Professional level

- Dedicated team
- Proactive
- Collecting and analysing incident data
- Identifying and escalating issues on organisational level
- Unifying and improving processes
- Transparency

# SRE @ Infobip

# Numbers

### Company

5 Business Areas

26 Requirement Areas

100+ teams

~900 engineers

~3000 employees total

### Products / Platform

23 products and channels

39 DCs

50+ locations

3 clouds (on prem + 2 public)

### Rate of change (monthly)

~30k deployments

~30B client interactions

~30k active web users

50+ maintenances

1 release (of all products)

# SRE

5 team members

30+ years of IB experience

10+ different IB job titles

50+ years of IT experience

# Platform monitoring

Platform, high-level alerts

Open channels to support and teams

Request teams to expose relevant metrics

Driving SLO adoption

# Incident Management

Owners of the IM process

Helping/Handling incidents

Incident commanders for complex incidents

Collection and analysis of incident meta-data

Monthly, quarterly, yearly reports

Post-incident reviews

# Tooling

Automating operational tasks

Automating processes

Reviewing usage of observability, alerting and escalation tools

Educating how to best use the above

# Coordination

High-impact and critical maintenances

Handling reliability-related inquiries by clients

Client integration when high levels of reliability is required

Mediators between stakeholders

# Culture

Blameless incident culture

Data-driven decision making

Contextualizing SRE practices for our way of work

Defining best practices: HA, monitoring, availability, etc.

Promoting a client-centric view of problems on our platform

# Culture

... if you do [have problems]

we will look for you

we will find you

and we will help you

# Product

Reliability review

Driving reliability improvements

# 1 godina SREće (1 year of happiness)

- 0 to 5 SREs
- Incident number increase
    - Improved detection
    - Improved reports
- On average, incident duration is halved when SRE member involved
    - No better metric ATM :(
- Reporting speed-up
    - Monthly: from weeks to < 1 day
    - Self-service dashboards for managers
- High, cross-company, visibility

# Road to SRE

# WHY SRE?

Fulfill contractual obligations

Proactively build and maintain reliable services
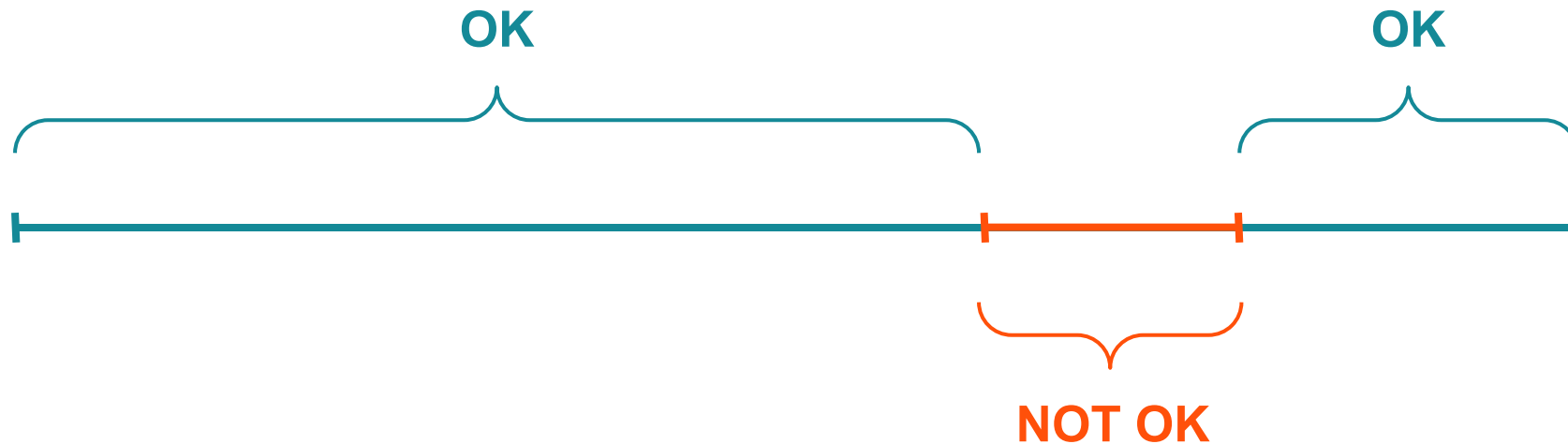
# WHY SRE?

hope != strategy

# WHAT IS NEEDED?



Reliability measures the functioning of a service over a period of time, under specified conditions.

# Measure

- Monitoring
- Observability is the base
- It is not trivial to have good measures
- If you don't have it, start today
- Once you have it, improve it constantly

*Progressive improvement beats delayed perfection.*
*Mark Twain*

# Function

- Categorization ok / not ok
- Define it carefully
- Make it observable
- Define it from the client PoV
  - Whoever or whatever the client in your context is

*If the client receives a 200 OK response, because the request was put into a queue, are they getting the service they paid for?*

# Incident

- Categorize incident / not incident
- When should people report and escalate problems?
- What are the thresholds?
- How to define them?
- They will happen
- Be ready
- SLOs make it easier to answer these

*Success is not final, failure is not fatal:*
*it is the courage to continue that counts.*
*Winston Churchill*

# Incident management

- Are all incidents equal?
- Are they equally severe?
- Do they have the same priority?
- How does one respond to an incident?
- Define the incident response process
  - Teach it
  - Exercise it
  - Improve it
- Guides
  - ITIL, ITSM, OODA
- Incident reports
  - For transparency
  - For improvements rather focus on specific incidents

# Disaster scenarios

- Disaster will happen too
- Are you ready?
    - Facebook was
- What is the cost of the service being down
    - Day?
    - Week?
    - Month?

*I think the environmental impact of this disaster is likely to have been very, very modest.*
*Tony Hayward, BP CEO*

# On-call

- Organizing incident response
- Protecting people
    - Their well-being
    - Their work-life balance
    - Ther happiness
- Organizing rotations
- Clear responsibilities
- Clear expectations

*Have you tried turning it off and on again?*
*IT Crowd*

# Culture

*Culture* eats *strategy* for breakfast.

Peter Drucker

# HOW TO SRE?

Each company does it differently

Needs to be aligned with the company culture

Start by adopting practices, one by one

# Post-mortems / Post-incident reviews

- Incidents are complex
- Incidents are unique
- Focus on finding all the contributing causes
    - There's rarely a single, root cause
- Define planned actions
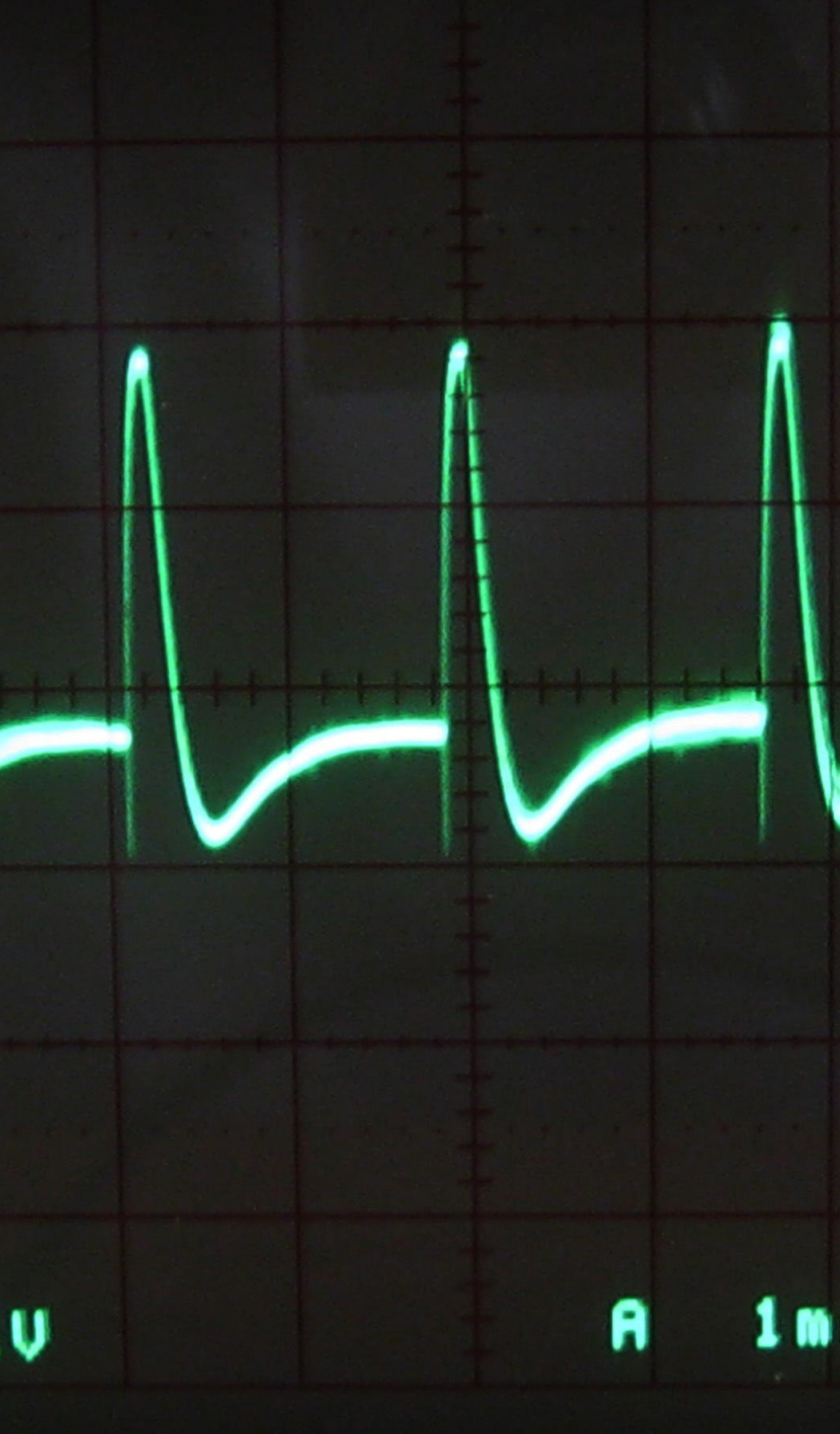    - Make sure they are executed

# SLO - Service Level Objective

- SLIs are not just another metric
    - All stakeholders agree on the importance
    - Relates to business value
- Does not have to be 100% precise
    - Constantly improve
- If objective is not met, actions are taken
- Set realistic targets
    - No point in failing, constantly
    - SLO is always stricter than SLA
- Do not make it a KPI / OKR

When a measure becomes a target,
it ceases to be a good measure.

Goodhart's Law

# Automated alerting

- Good alerting is hard to set up
- Requires maintenance
- Requires constant improvements
- To properly scale, requires a strategy
- Differentiate between
    - Alerts - as few as possible
    - Troubleshooting data - as much as CEO is willing to pay
    - Notifications - calling you in the middle of the night
- Not all metrics need to be alerted on
- Not all alerts need to trigger notifications
- Not all alerts need to trigger end-of-the-world notifications

*Be alert… the world needs more lerts.*
*Woody Allen*

# Chaos engineering

- Handy practice
- Uncovers some types of problems
- Can be used to improve reliability
- Start practicing when you think you are reliable enough

# Other practices

- Fire drills, game days, disaster recovery tests
- MTTx metrics
    - Collection and analysis
    - Beware of averages
    - TTx histograms
- Service and organisation registry
    - Keep the two in sync
- Eliminating toil
    - Meta-practice, should be included in everything
- Data analysis, statistics
- Many more

*Average: a random number that falls somewhere between the maximum and 1/2 the median. Most often used to ignore reality.*
*Gil Tene*

k8s?

The Google Model
We Are Now SRE
SRE Center of Practice / Excellence
Embedded SRE
[Github repo: How They SRE](...)...
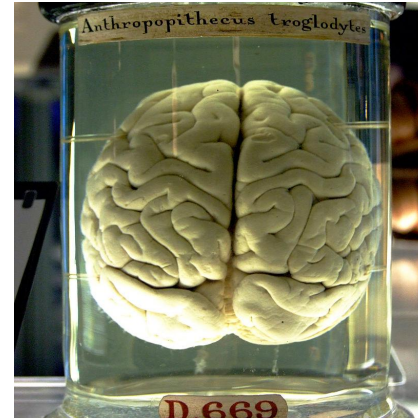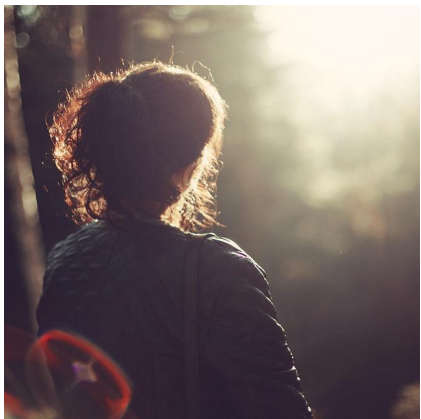
# Epilogue

**Start small**

- Start with anything
- Start with what you have
- Practice makes perfect
- Improve constantly

**Start smart**

- Do not reinvent the wheel
- There's tons of resources
- Do not ignore the history of how practices evolved and why

**Be kind to yourself**

- SRE, not a role, a condition
- Handling production is challenging
- Lots of context switching
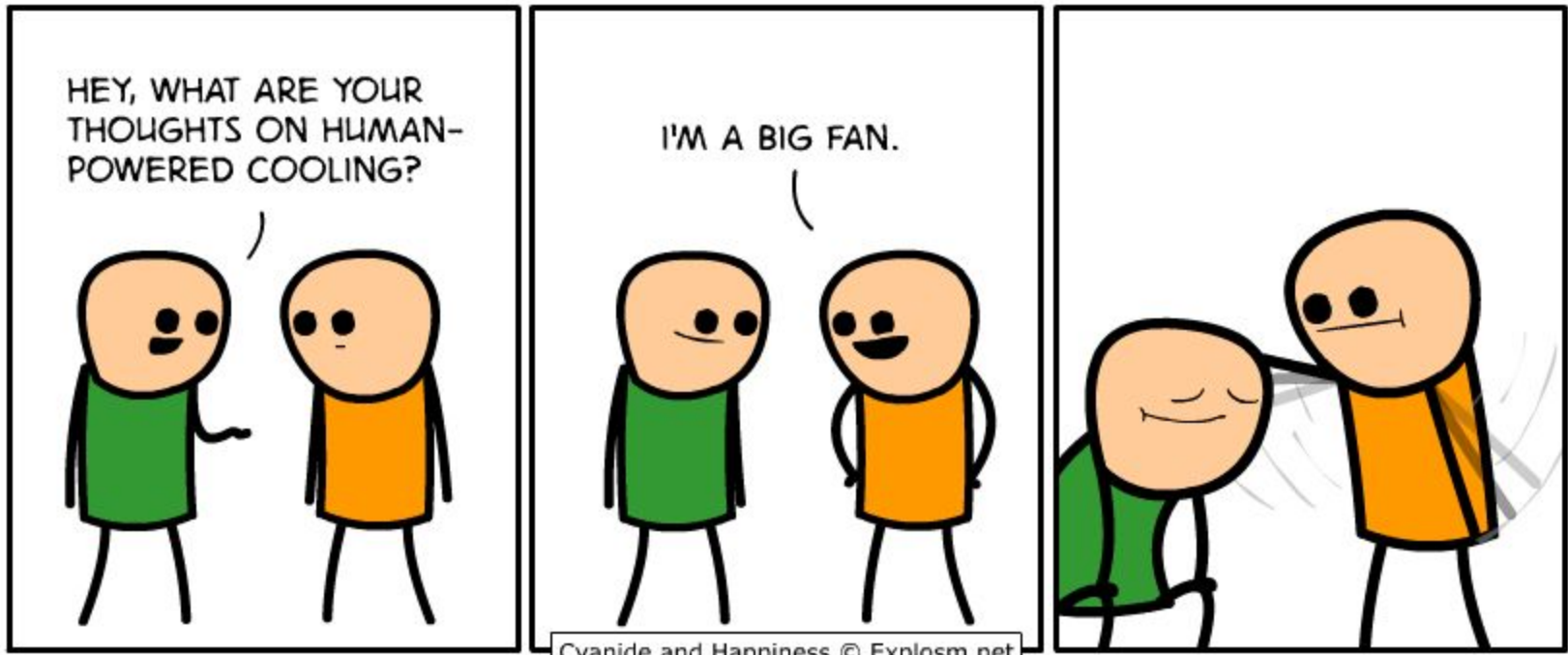- Lots of quick decisions
- Can be highly stressful

**Be kind to others**

- Communication is essential
- Lots of stakeholders and their specific dialects
- Lots of different cultures
- Make people responsible for their actions and services

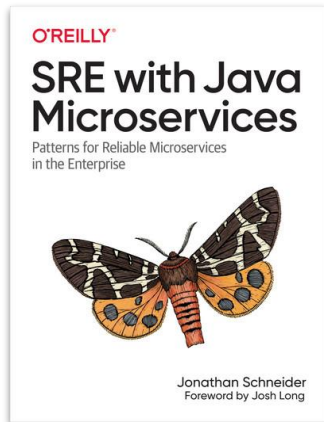Strategy for engaging humans doing ops with something worthy of their mental capacity
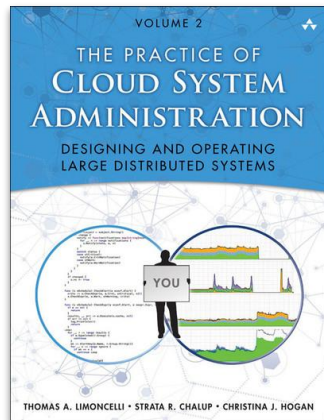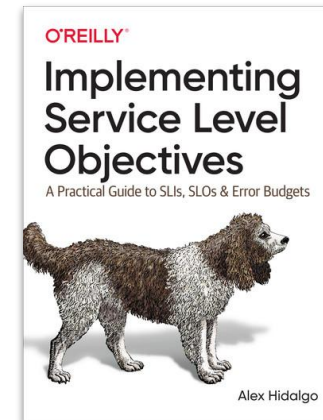
# THANK YOU!

infobip

# References



- Code-level / microservice architecture reliability
- Observability: deep-dive
- Recommended for all SW engineers



- Stability patterns
- Examples of real-life failures and how to mitigate them
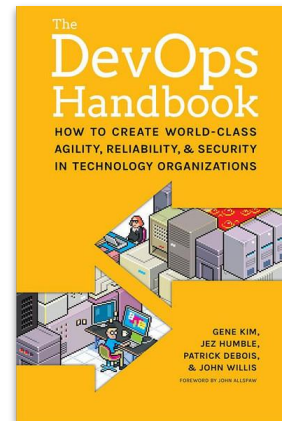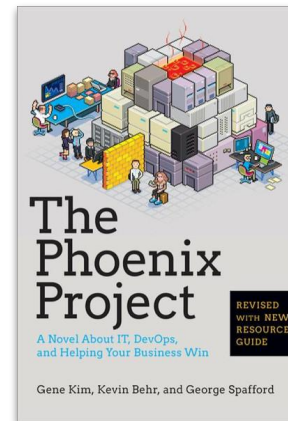- Recommended for SW engineers working with distributed systems



- Overview of ops required in the cloud
- Design, operate, assess, improve
- Recommended for tech-savvy managers, new / evolving Sys Amins, SREs, Devs doing Ops
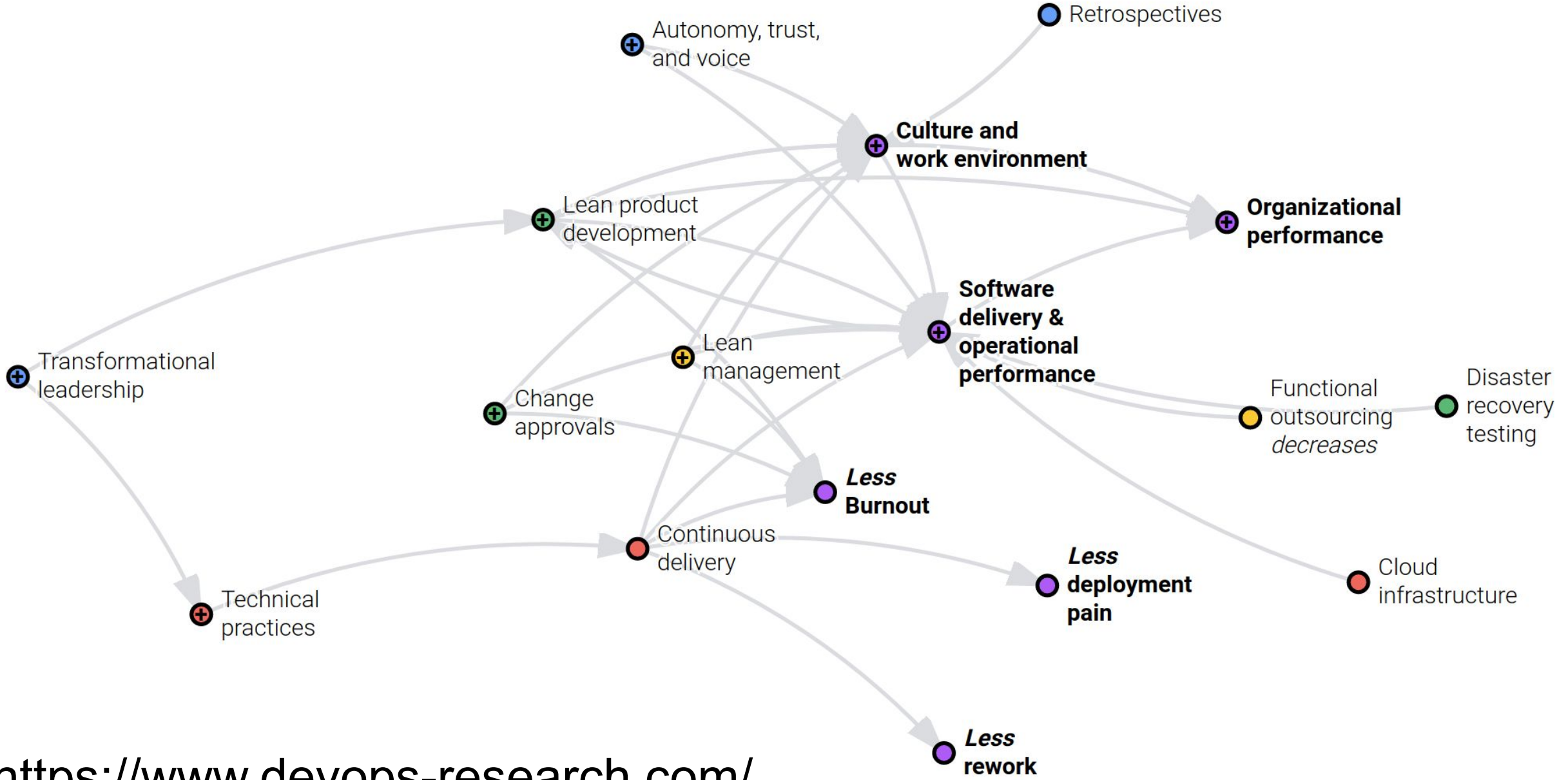


- All you need to know about SLOs
- Deep-dive into the subject
- Recommended for managers, architects and senior SW engineers

Newsletter: https://sreweekly.com/

- Phoenix project: DevOps explained as a fictional story
- Handbook: Why, what and how to DevOps
- Recommended to all thinking they need to hire a DevOp

Autonomy, trust, and voice

Retrospectives

Culture and work environment

Organizational performance

Lean product development

Software delivery & operational performance

Transformational leadership

Lean management

Functional outsourcing *decreases*

Disaster recovery testing

Change approvals

*Less* Burnout

Continuous delivery

*Less* deployment pain

Cloud infrastructure

Technical practices

*Less* rework

https://www.devops-research.com/